

REGRESSION PART 5: BIAS-VARIANCE TRAIDOFF

Hsin-Min Lu

盧信銘

台大資管系

The Squared Loss Function (Section 1.5.5)

- Using RMSE or MSE as performance measure meaning that we are looking at squared loss function.
- Consider a scenario when we want to study the relationship between feature vector x and outcome t .
- We can think of x and t as random variables when we are collecting data points. They have a joint distribution $p(x, t)$.
- We want to construct our prediction function $y(x)$ by minimizing the expected squared loss:
- $$E[L] = \int \int (y(x) - t)^2 p(x, t) dt dx$$
- **Meaning: minimize prediction error for all data.**



The Squared Loss Function

- $E[L] = \iint (y(x) - t)^2 p(x, t) dt dx$
$$= \iint (y(x) - t)^2 p(t|x) p(x) dt dx$$
- Decompose $(y(x) - t)^2 = (y(x) - E[t|x] + E[t|x] - t)^2$
- $= (y(x) - E[t|x])^2 + 2(y(x) - E[t|x])(E[t|x] - t) + (E[t|x] - t)^2$
- Taking expectation on the three terms separately
- $E[(y(x) - E[t|x])^2] = \int \int (y(x) - E[t|x])^2 p(t|x) p(x) dt dx$
 $= \int (y(x) - E[t|x])^2 p(x) dx$ (why???)



The Squared Loss Function

- Recall $(y(x) - t)^2 = (y(x) - E[t|x])^2 + 2(y(x) - E[t|x])(E[t|x] - t) + (E[t|x] - t)^2$
- $E[2(y(x) - E[t|x])(E[t|x] - t)]$
$$= \int \int 2(y(x) - E[t|x])(E[t|x] - t) p(t|x) p(x) dt dx$$
$$= \int 2(y(x) - E[t|x])(E[t|x] - E[t|x]) p(x) dx = 0$$
- $E[(E[t|x] - t)^2]$
$$= \iint (E[t|x] - t)^2 p(t|x) p(x) dt dx = \int \text{Var}[t|x] p(x) dx$$



The Squared Loss Function

- Putting all three terms together
- $E[L] = \int (y(x) - E[t|x])^2 p(x) dx + 0 + \int \text{Var}[t|x] p(x) dx$
- The last term is the average of noise across all possible x .
 - We can do nothing about this term.
- We can minimize $E[L]$ by setting $y(x)$ to $E[t|x]$.
- That is, the best prediction that minimize the expected squared loss is $y(x) = E[t|x]$.
- If this is the case, then the expected square loss, or the expected MSE, is $\int \text{Var}[t|x] p(x) dx$



3.2 The Bias-Variance Decomposition (1)

- Recall the *expected squared loss*,
- $E[L] = \int (y(x) - E[t|x])^2 p(x) dx + \int \text{Var}[t|x] p(x) dx$
- The second term of $E[L]$ corresponds to the noise (variance) inherent in the random variable t .
independent of the choice of $y(x)$
- Can we set $y(x) = E[t|x]$?
- In theory, yes, but...
- In reality, we do not know $E[t|x]$ for sure.



The Bias-Variance Decomposition (2)

- In reality, we are given limited dataset in order to learn $E[t|x]$.
- Following textbook notation, let $h(x) = E[t|x]$
- Suppose we were given multiple data sets, each of size N . Any particular data set, \mathcal{D} , will give a particular function $y(x; \mathcal{D})$. We then have

$$\begin{aligned} & \{y(\mathbf{x}; \mathcal{D}) - h(\mathbf{x})\}^2 \\ &= \{y(\mathbf{x}; \mathcal{D}) - \mathbb{E}_{\mathcal{D}}[y(\mathbf{x}; \mathcal{D})] + \mathbb{E}_{\mathcal{D}}[y(\mathbf{x}; \mathcal{D})] - h(\mathbf{x})\}^2 \\ &= \{y(\mathbf{x}; \mathcal{D}) - \mathbb{E}_{\mathcal{D}}[y(\mathbf{x}; \mathcal{D})]\}^2 + \{\mathbb{E}_{\mathcal{D}}[y(\mathbf{x}; \mathcal{D})] - h(\mathbf{x})\}^2 \\ &\quad + 2\{y(\mathbf{x}; \mathcal{D}) - \mathbb{E}_{\mathcal{D}}[y(\mathbf{x}; \mathcal{D})]\}\{\mathbb{E}_{\mathcal{D}}[y(\mathbf{x}; \mathcal{D})] - h(\mathbf{x})\}. \end{aligned}$$

Cross product term



The Bias-Variance Decomposition (3)

- Taking the expectation over D , the cross product term vanished:

- $$E_D [2\{y(x; D) - E_D[y(x; D)]\}\{E_D[y(x; D)] - h(x)\}] = 2[E_D[y(x; D)] - E_D[y(x; D)]]\{E_D[y(x; D)] - h(x)\} = 0$$

- Thus we have:

$$\begin{aligned} & \mathbb{E}_D [\{y(\mathbf{x}; \mathcal{D}) - h(\mathbf{x})\}^2] \\ &= \underbrace{\{\mathbb{E}_D[y(\mathbf{x}; \mathcal{D})] - h(\mathbf{x})\}^2}_{(\text{bias})^2} + \underbrace{\mathbb{E}_D [\{y(\mathbf{x}; \mathcal{D}) - \mathbb{E}_D[y(\mathbf{x}; \mathcal{D})]\}^2]}_{\text{variance}}. \end{aligned}$$



The Bias-Variance Decomposition (4)

- Putting everything together, we can write

$$\text{expected loss} = (\text{bias})^2 + \text{variance} + \text{noise}$$

- where

$$h(x) = E[t|x]$$

$$(\text{bias})^2 = \int \{\mathbb{E}_{\mathcal{D}}[y(\mathbf{x}; \mathcal{D})] - h(\mathbf{x})\}^2 p(\mathbf{x}) d\mathbf{x}$$

$$\text{variance} = \int \mathbb{E}_{\mathcal{D}} [\{y(\mathbf{x}; \mathcal{D}) - \mathbb{E}_{\mathcal{D}}[y(\mathbf{x}; \mathcal{D})]\}^2] p(\mathbf{x}) d\mathbf{x}$$

$$\text{noise} = \iint \{h(\mathbf{x}) - t\}^2 p(\mathbf{x}, t) d\mathbf{x} dt$$

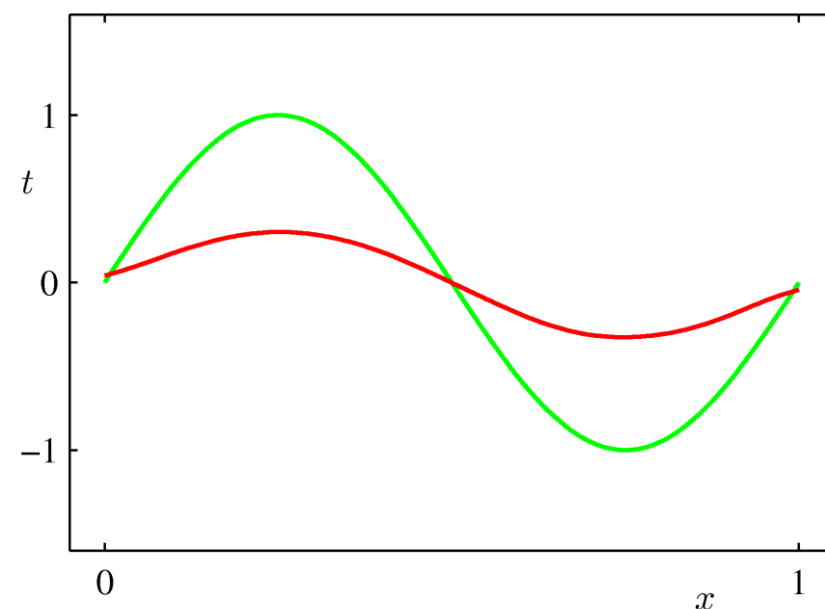
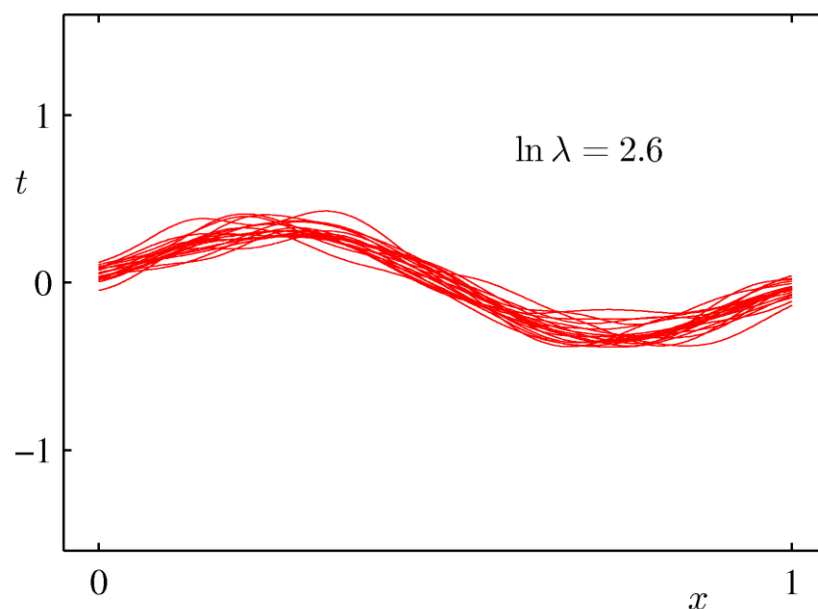


The Bias-Variance Decomposition (5a)

- Example: 100 data sets (each having $N=25$ data points) from the sinusoidal, varying the degree of regularization, λ . $M=25$, (24 Gaussian basis functions).
- Fitting data via minimizing $\frac{1}{2} \sum (t_n - \mathbf{w}^T \phi(\mathbf{x}_n))^2 + \frac{\lambda}{2} \sum w_j^2$
- Higher $\lambda \rightarrow$ More rigid model

The Bias-Variance Decomposition (5b)

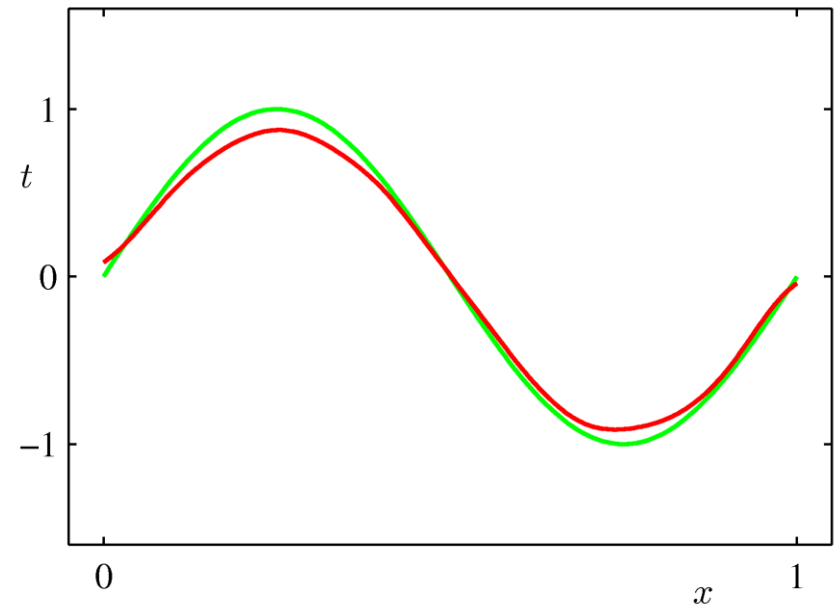
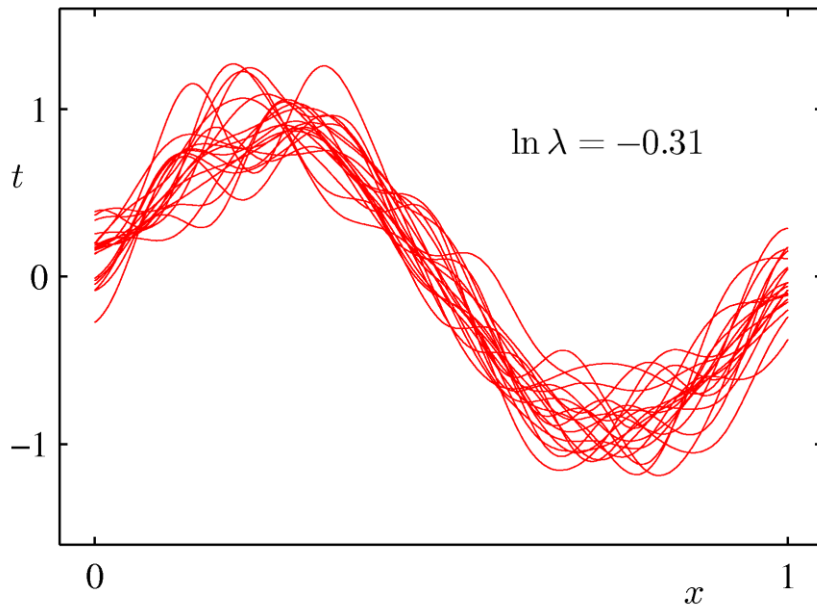
- $\ln \lambda = 2.6$ (low variance, high bias)



Green: true model
Red: average of 100 fits

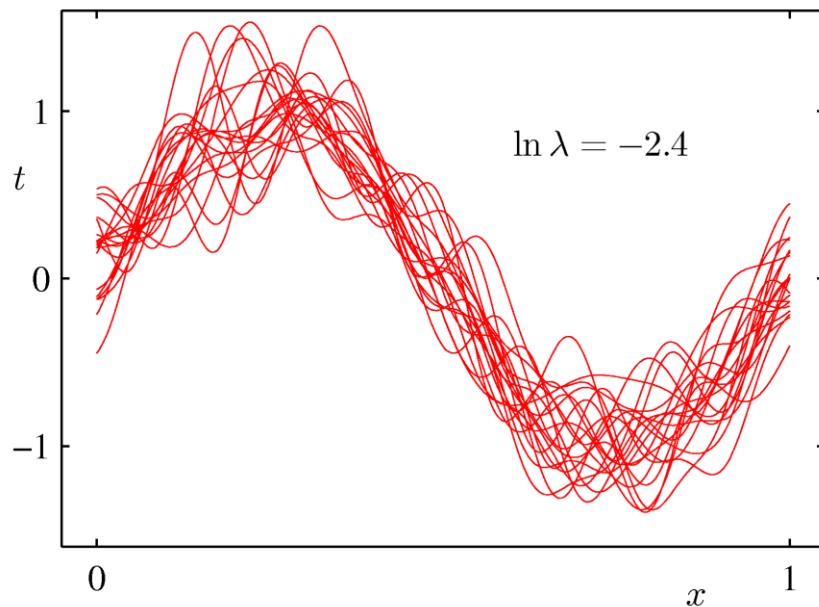
The Bias-Variance Decomposition (6)

- $\ln \lambda = -0.31$ (smaller regularization)
- Higher variance, lower bias



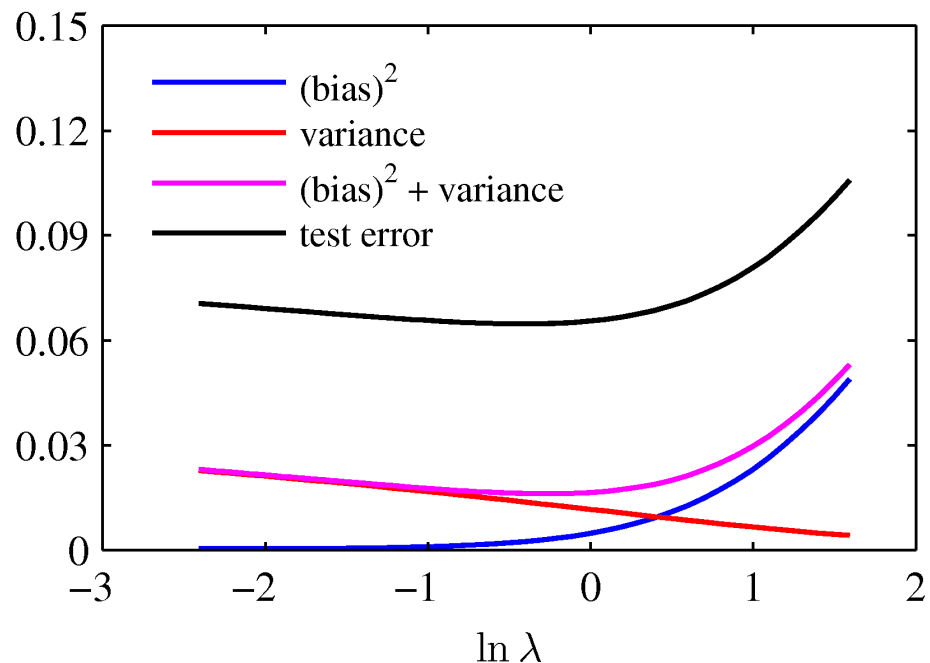
The Bias-Variance Decomposition (7)

- $\ln \lambda = -2.4$ (even smaller regularization)
- Even higher variance, even lower bias



The Bias-Variance Trade-off

- From these plots, we note that an over-regularized model (large λ) will have a high bias, while an under-regularized model (small λ) will have a high variance.
- We can usually gain some level of prediction accuracy by let-go unbiasedness and select a reasonable level of bias-variance trade-off.



Reading List

- PRML Ch 1.5.5, Ch 3.2